

CLAIMS

What is claimed is:

1. A method for processing a formatted computer-readable source document having one or more pages, the source document having text comprising a plurality of words, each word comprising one or more characters, each character having a character appearance defined by one or more font properties and each word having a word appearance defined by the font properties of its characters and a position of the word on a page, the method comprising:
- partitioning the formatted text into one or more groups of words and assigning at least one element from a predefined set of markup language elements to each group of words, each group of words having a group appearance defined by one or more text properties;
 - for each of the predefined elements that is assigned to at least one group or words, deriving an element style comprising a character style, a layout style or both, the character style being derived from the font properties of the characters of the words in the groups of words to which the element is assigned, and the layout style being derived from the text properties of the groups of words to which the element is assigned; and
 - creating an electronic document comprising a style sheet defining each of the element styles.
2. The method of claim 1, wherein the text is partitioned into groups of words according to word positions and an element is assigned to each group of words solely on the basis of the position of the group of words on the page.
3. The method of claim 1, wherein the text is partitioned into groups of words according to the font properties of the words in the text and an element is assigned to each group of words solely on the basis of the font properties of the words in the group of words.
4. The method of claim 1, wherein the text is partitioned into groups of words according to word positions and the font properties of the words in the text and at least one element is assigned to each group of words based on the position of the group of words on the page and the font properties of the words in the group of words.

66001-1109450

- 1 5. The method of claim 1, further comprising:
2 for each element assigned to at least one group of words, comparing the group
3 appearances of all groups of words to which the element is assigned, creating one or more
4 alternate elements if the differences among the group appearances exceed a predefined
5 threshold, and assigning each group of words to the original element or an alternate element.
- 1 6. The method of claim 5, wherein a numeric value defining the predefined threshold is
2 obtained from user input.
- 1 7. The method of claim 5, wherein a numeric value defining the predefined threshold is a
2 preprogrammed numeric value.
- 1 8. The method of claim 1, wherein the set of predefined elements is the set of HyperText
2 Markup Language elements defined in HTML 4.0.
- 1 9. The method of claim 1, wherein the set of predefined elements is the set of Extensible
2 Markup Language elements defined in XML 1.0.
- 1 10. The method of claim 1, wherein the predefined set of elements comprises:
2 a header element and a paragraph element.
- 1 11. The method of claim 10, wherein the predefined set of elements further comprises:
2 an address element, a blockquote element, a list element, a table element and a
3 caption element.
- 1 12. The method of claim 1, wherein the character style comprises at least one font property
2 and an associated value.

- 1 13. The method of claim 12, wherein the font property is selected from a predefined set of
2 font properties comprising a font family, a font style, a font weight, a font variant and a font
3 size.
- 1 14. The method of claim 1, wherein the layout style comprises at least one text property and
2 an associated value.
- 1 15. The method of claim 14, wherein the text property is selected from a predefined set of
2 text properties comprising a text decoration, a text alignment, a text indentation and a text
3 transformation as defined in CSS1.
- 4 16. The method of claim 1, wherein the source document is derived from a raster image of
5 the text processed by an optical character recognition (OCR) process.
- 1 17. The method of claim 16, further comprising translating the raster image of the text into a
2 HTML description of the text.
- 1 18. The method of claim 1, further comprising detecting and setting page margins for the
2 document page.
- 1 19. The method of claim 1, wherein the style sheet is an extensible style sheet (XSL).
- 1 20. The method of claim 1, further comprising creating an electronic document comprising a
2 markup language version of the source document.
- 1 21. A computer program product, tangibly stored on a computer-readable medium, for
2 processing a formatted computer-readable source document having one or more pages, the
3 source document having text comprising a plurality of words, each word comprising one or
4 more characters, each character having a character appearance defined by one or more font
5 properties and each word having a word appearance defined by the font properties of its

6 characters and a position of the word on a page, the product comprising instructions operable
7 to cause a programmable processor to:

8 partition the formatted text into one or more groups of words and assign at least one
9 element from a predefined set of markup language elements to each group of words, each
10 group of words having a group appearance defined by one or more text properties;

11 for each of the predefined elements that is assigned to at least one group of words,
12 derive an element style comprising a character style, a layout style or both, the character
13 style being derived from the font properties of the characters of the words in the groups of
14 words to which the element is assigned, and the layout style being derived from the text
15 properties of the groups of words to which the element is assigned; and

16 create an electronic document comprising a style sheet defining each of the element
17 styles.

1 22. The product of claim 21, wherein the text is partitioned into groups of words according
2 to word positions and an element is assigned to each group of words solely on the basis of
3 the position of the group of words on the page.

1 23. The product of claim 21, wherein the text is partitioned into groups of words according
2 to the font properties of the words in the text and an element is assigned to each group of
3 words solely on the basis of the font properties of the words in the group of words.

1 24. The product of claim 21, wherein the text is partitioned into groups of words according
2 to word positions and the font properties of the words in the text and at least one element is
3 assigned to each group of words based on the position of the group of words on the page and
4 the font properties of the words in the group of words.

1 25. The product of claim 21, further comprising:

2 for each element assigned to at least one group of words, comparing the group
3 appearances of all groups of words to which the element is assigned, creating one or more
4 alternate elements if the differences among the group appearances exceed a predefined
5 threshold, and assigning each group of words to the original element or an alternate element.